

第十章

极大似然估计

Outline

- 1 Maximum Likelihood Estimation
- 2 Conditional ML Estimation of the Classical Linear Regression Model
- 3 Score, Hessian, and Information
- 4 Consistency and Asymptotic Normality of MLE
- 5 Wald, LM, and LR Test

Maximum Likelihood Estimation

The basic idea of the ML principle is to choose the parameter estimates to maximize the probability of obtaining the data.

定义 (Model)

A model for a random sample is the assumption that X_i , $i = 1, \dots, n$, are i.i.d. with known density function $f(x|\theta)$ or mass function $\pi(x|\theta)$ with unknown parameter $\theta \in \Theta \subseteq \mathbb{R}^p$.

定义 (Correctly Specified Model)

A model is **correctly specified** when there is a unique parameter value $\theta_0 \in \Theta$ such that $f(x|\theta_0) = f(x)$, the true data distribution. This parameter value θ_0 is called the true parameter value. The parameter θ_0 is **unique** if there is no other θ such that $f(x|\theta_0) = f(x|\theta)$. A model is **mis-specified** if there is no parameter value $\theta \in \Theta$ such that $f(x|\theta) = f(x)$.

Likelihood

The joint density evaluated at the observed data (x_1, x_2, \dots, x_n) and viewed as a function of θ is called the likelihood function:

$$L_n(\theta) \equiv f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

for continuous random variables and

$$L_n(\theta) \equiv \prod_{i=1}^n \pi(x_i | \theta)$$

for discrete random variables.

The value of θ most compatible with the observations is the value which maximizes the likelihood.

定义 (MLE)

The maximum likelihood estimator $\hat{\theta}$ of θ is the value which maximizes $L_n(\theta)$:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} L_n(\theta). \quad (1)$$

In most cases it is more convenient to calculate and maximize the logarithm of the likelihood function rather than the level of the likelihood.

定义 (Log-likelihood)

The log-likelihood function is

$$\ell_n(\theta) \equiv \log L_n(\theta) = \sum_{i=1}^n \log f(x_i | \theta).$$

The maximizer of the likelihood and log-likelihood function are the same, since the logarithm is a monotonically increasing function.

Define the expected log density function

$$\ell(\boldsymbol{\theta}) = \mathbb{E}[\log f(X|\boldsymbol{\theta})].$$

Theorem

When the model is correctly specified the true parameter $\boldsymbol{\theta}_0$ maximizes the expected log density $\ell(\boldsymbol{\theta})$:

$$\boldsymbol{\theta}_0 = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \ell(\boldsymbol{\theta}).$$

The sample analog of $\ell(\boldsymbol{\theta})$ is the average log-likelihood

$$\bar{\ell}_n(\boldsymbol{\theta}) = \frac{1}{n} \ell_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \log f(x_i|\boldsymbol{\theta})$$

which has maximizer $\hat{\boldsymbol{\theta}}$.

证明.

Take any $\theta \neq \theta_0$. Since the difference of logs is the log of the ratio

$$\begin{aligned}\ell(\theta) - \ell(\theta_0) &= \mathbb{E} [\log(f(X|\theta)) - \log(f(X|\theta_0))] \\ &= \mathbb{E} \left[\log \left(\frac{f(X|\theta)}{f(X|\theta_0)} \right) \right] < \log \left(\mathbb{E} \left[\frac{f(X|\theta)}{f(X|\theta_0)} \right] \right),\end{aligned}$$

where the last inequality follows from Jensen's inequality and is strict since $f(x|\theta) \neq f(x|\theta_0)$ with positive probability. Note that

$$\begin{aligned}\log \left(\int \frac{f(x|\theta)}{f(x|\theta_0)} f(x) dx \right) &= \log \left(\int \frac{f(x|\theta)}{f(x|\theta_0)} f(x|\theta_0) dx \right) \\ &= \log \left(\int f(x|\theta) dx \right) = 0.\end{aligned}$$

Therefore, $\ell(\theta) < \ell(\theta_0)$. □

Invariance Property

A special property of the MLE is that it is invariant to transformations.

Theorem

If $\hat{\theta}$ is the MLE of θ then for any transformation $\beta = h(\theta)$ the MLE of β is $\hat{\beta} = h(\hat{\theta})$.

证明.

We can write the likelihood for the transformed parameter as

$$L_n^*(\beta) = \max_{h(\theta)=\beta} L_n(\theta).$$

The MLE for β maximizes $L_n^*(\beta)$. Evaluating $L_n^*(\beta)$ at $h(\hat{\theta})$ we find

$$L_n^*(h(\hat{\theta})) = \max_{h(\theta)=h(\hat{\theta})} L_n(\theta) = L_n(\hat{\theta}).$$

The final equality holds because $\theta = \hat{\theta}$ satisfies $h(\theta) = h(\hat{\theta})$ and maximizes $L_n(\theta)$. □

To find the MLE we take the following steps:

- 1 Construct $f(x|\theta)$ as a function of x and θ
- 2 Take the logarithm $\log f(x|\theta)$
- 3 Evaluate at $x = x_i$ and sum over i : $\ell_n(\theta) = \sum_{i=1}^n \log f(x_i|\theta)$
- 4 If possible, solve the first order condition (F.O.C.) to find the maximum
- 5 If solving the F.O.C. is not possible use numerical methods to maximize $\ell_n(\theta)$

Conditional ML Estimation of the Classical Linear Regression Model

Conditional Maximum Likelihood

In most applications, observation for i is partitioned into two groups, y_i and \mathbf{x}_i . Let $f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)$ be the conditional density of y_i , given \mathbf{x}_i , and let $f(\mathbf{x}_i; \boldsymbol{\psi}_0)$ be the marginal density of \mathbf{x}_i . Then

$$f(y_i, \mathbf{x}_i; \boldsymbol{\theta}_0, \boldsymbol{\psi}_0) = f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)f(\mathbf{x}_i; \boldsymbol{\psi}_0).$$

Let $\mathbf{w}_i = (y_i, \mathbf{x}_i)$. The log likelihood is

$$\sum_{i=1}^n \log f(\mathbf{w}_i; \boldsymbol{\theta}, \boldsymbol{\psi}) = \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) + \sum_{i=1}^n \log f(\mathbf{x}_i; \boldsymbol{\psi}).$$

The first term on the right-hand side is the log conditional likelihood. **The conditional ML estimator of $\boldsymbol{\theta}_0$ maximizes this first term.** If the second term does not depend on $\boldsymbol{\theta}$, the conditional ML estimator of $\boldsymbol{\theta}_0$ is numerically the same as the joint ML estimator.

Under the assumptions of classical linear regression, we have

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta}_0 + \boldsymbol{\varepsilon} \\ \mathbf{y} \mid \mathbf{X} &\sim N(\mathbf{X}\boldsymbol{\beta}_0, \sigma_0^2 \mathbf{I}_n). \end{aligned}$$

Thus, the conditional density of \mathbf{y} given \mathbf{X} is¹

$$f(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\beta}_0, \sigma_0^2) = (2\pi\sigma_0^2)^{-n/2} \exp \left[-\frac{1}{2\sigma_0^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0) \right].$$

Therefore, the log likelihood function is

$$\log L(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

¹The density function for an n -variate normal distribution $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is

$$(2\pi)^{-n/2} |\boldsymbol{\Sigma}|^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right].$$

ML via Concentrated Likelihood

The log likelihood function can be maximized in two stages:

- ① maximize over β for any given σ^2 ;
- ② maximize over σ^2 taking into account that the β obtained in the first stage could depend on σ^2 .

The log likelihood function in which β is constrained to be the value from the first stage is called **the concentrated log likelihood function**. It is easy to see that

$$\text{concentrated log likelihood} = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} e'e.$$

定理 (ML Estimator of (β, σ^2))

Suppose Assumptions 1.1–1.5 hold. Then $\hat{\beta}_{MLE} = \mathbf{b}$ and $\hat{\sigma}_{MLE}^2 = e'e/n$.

Score, Hessian, and Information

The **likelihood score** is the derivative of the likelihood function:

$$\mathbf{S}_n(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \ell_n(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \log f(x_i | \boldsymbol{\theta}).$$

Note that $\mathbf{S}_n(\hat{\boldsymbol{\theta}}) = \mathbf{0}$ when $\hat{\boldsymbol{\theta}}$ is an interior solution.

The **likelihood Hessian** is the negative second derivative of the likelihood function:

$$\mathcal{H}_n(\boldsymbol{\theta}) = -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \ell_n(\boldsymbol{\theta}) = -\sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log f(x_i | \boldsymbol{\theta}).$$

The **efficient score** is the derivative of the log-likelihood for a **single observation**, evaluated at the random variable X and the true parameter vector:

$$\mathbf{S} = \frac{\partial}{\partial \boldsymbol{\theta}} \log f(X | \boldsymbol{\theta}_0).$$

Theorem

Assume that the model is correctly specified, the support of X does not depend on θ , and θ_0 lies in the interior of Θ . Then the efficient score S satisfies $\mathbb{E}[S] = \mathbf{0}$.

证明.

$$\mathbb{E}[S] = \mathbb{E} \left[\frac{\partial}{\partial \theta} \log f(X|\theta_0) \right] = \frac{\partial}{\partial \theta} \mathbb{E} [\log f(X|\theta_0)] = \frac{\partial}{\partial \theta} \ell(\theta_0) = \mathbf{0}.$$



- The assumption that the support does not depend on the parameter is needed for interchange of integration and differentiation.
- The assumption that θ_0 lies in the interior of Θ is needed to ensure that the expected log-likelihood $\ell(\theta)$ satisfies a first-order-condition.

The **Fisher information** is the variance of the efficient score:

$$\mathcal{I}_{\theta} = \mathbb{E} [SS'] .$$

The **expected Hessian** is

$$\mathcal{H}_{\theta} = -\frac{\partial^2}{\partial \theta \partial \theta'} \ell(\theta_0) .$$

Under regularity conditions,

$$\mathcal{H}_{\theta} = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta \partial \theta'} \log f(X|\theta_0) \right] .$$

Theorem: Information Matrix Equality

Assume that the model is correctly specified and the support of X does not depend on θ . Then the Fisher information equals the expected Hessian:

$$\mathcal{I}_{\theta} = \mathcal{H}_{\theta}.$$

证明.

The expected Hessian equals

$$\begin{aligned}\mathcal{H}_{\theta} &= -\frac{\partial^2}{\partial \theta \partial \theta'} \mathbb{E} [\log f(X|\theta_0)] = -\frac{\partial}{\partial \theta} \mathbb{E} \left[\frac{\frac{\partial}{\partial \theta'} f(X|\theta)}{f(X|\theta)} \right] \Big|_{\theta=\theta_0} \\ &= \underbrace{-\frac{\partial}{\partial \theta} \mathbb{E} \left[\frac{\frac{\partial}{\partial \theta'} f(X|\theta)}{f(X|\theta_0)} \right] \Big|_{\theta=\theta_0}}_0 + \underbrace{\mathbb{E} \left[\frac{\frac{\partial}{\partial \theta} f(X|\theta_0) \frac{\partial}{\partial \theta'} f(X|\theta_0)}{f(X|\theta_0)^2} \right]}_{\mathcal{I}_{\theta}}.\end{aligned}$$



Cramér-Rao Lower Bound

The information matrix provides a lower bound for the variance among unbiased estimators.

Theorem: Cramér-Rao Lower Bound

Assume that the model is correctly specified, the support of X does not depend on θ , and θ_0 lies in the interior of Θ . If $\tilde{\theta}$ is an **unbiased estimator of θ** then $\text{var}[\tilde{\theta}] \geq (n\mathcal{I}_{\theta})^{-1}$.

- The Cramér-Rao Lower Bound (CRLB) is $(n\mathcal{I}_{\theta})^{-1}$.
- An estimator $\tilde{\theta}$ is Cramér-Rao efficient if it is unbiased for θ and $\text{var}[\tilde{\theta}] = (n\mathcal{I}_{\theta})^{-1}$.

证明.

Write $\mathbf{x} = (x_1, \dots, x_n)'$ and $\mathbf{X} = (X_1, \dots, X_n)'$. The joint density of \mathbf{X} is $f(\mathbf{x}|\boldsymbol{\theta}) = f(x_1|\boldsymbol{\theta}) \cdots f(x_n|\boldsymbol{\theta})$. The likelihood score is $\mathbf{S}_n(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}|\boldsymbol{\theta})$. Unbiasedness of $\tilde{\boldsymbol{\theta}}$ means that, for all $\boldsymbol{\theta}$,

$$\boldsymbol{\theta} = \mathbb{E}_{\boldsymbol{\theta}}[\tilde{\boldsymbol{\theta}}(\mathbf{X})] = \int \tilde{\boldsymbol{\theta}}(\mathbf{x}) f(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x}.$$

The vector derivative of the left side is

$$\frac{\partial}{\partial \boldsymbol{\theta}'} \boldsymbol{\theta} = \mathbf{I}_p.$$

and the vector derivative of the right side is

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}'} \int \tilde{\boldsymbol{\theta}}(\mathbf{x}) f(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} &= \int \tilde{\boldsymbol{\theta}}(\mathbf{x}) \frac{\partial}{\partial \boldsymbol{\theta}'} f(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} \\ &= \int \tilde{\boldsymbol{\theta}}(\mathbf{x}) \frac{\partial}{\partial \boldsymbol{\theta}'} \log f(\mathbf{x}|\boldsymbol{\theta}) f(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x}. \end{aligned}$$

证明.

Evaluated at the true value this is

$$\mathbb{E} \left[\tilde{\boldsymbol{\theta}}(\mathbf{X}) \mathbf{S}_n(\boldsymbol{\theta}_0)' \right] = \mathbb{E} \left[\left(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right) \mathbf{S}_n(\boldsymbol{\theta}_0)' \right], \quad \text{since } \mathbb{E} [\mathbf{S}_n(\boldsymbol{\theta}_0)] = \mathbf{0}.$$

Therefore, the variance matrix of stacked $\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0$ and $\mathbf{S}_n(\boldsymbol{\theta}_0)$ is $\begin{bmatrix} \mathbf{V} & \mathbf{I}_p \\ \mathbf{I}_p & n\mathcal{J}_{\boldsymbol{\theta}} \end{bmatrix}$, where $\mathbf{V} = \text{var}[\tilde{\boldsymbol{\theta}}]$. Since it is a variance matrix it is positive semi-definite and thus remains so if we pre- and post-multiply by the same matrix, e.g.

$$\begin{bmatrix} \mathbf{I}_p & -(n\mathcal{J}_{\boldsymbol{\theta}})^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{V} & \mathbf{I}_p \\ \mathbf{I}_p & n\mathcal{J}_{\boldsymbol{\theta}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_p \\ -(n\mathcal{J}_{\boldsymbol{\theta}})^{-1} \end{bmatrix} = \mathbf{V} - (n\mathcal{J}_{\boldsymbol{\theta}})^{-1} \geq \mathbf{0}.$$



Cramér-Rao Lower Bound for \mathbf{b}

定理 (\mathbf{b} is the Best Unbiased Estimator (BUE))

Under Assumptions 1.1–1.5, the OLS estimator \mathbf{b} is BUE in that any other unbiased (but not necessarily linear) estimator has larger conditional variance in the matrix sense.

The ML estimator of σ_0^2 is biased, so the Cramér-Rao bound does not apply.

Consistency and Asymptotic Normality of MLE

Uniform Convergence in Probability

定义

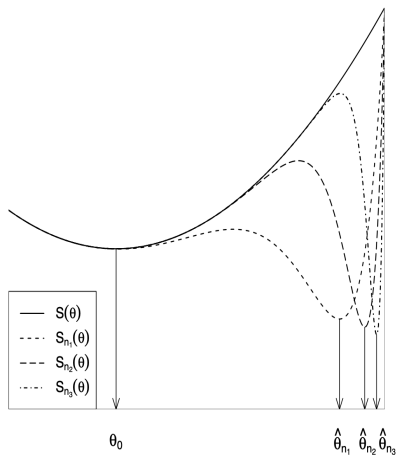
Pointwise convergence in probability of $\bar{\ell}_n(\cdot)$ to $\ell(\cdot)$ means that the sequence of random variables $|\bar{\ell}_n(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta})|$ ($n = 1, 2, \dots$) converges in probability to 0 for each $\boldsymbol{\theta}$. **Uniform convergence in probability** means

$$\sup_{\boldsymbol{\theta} \in \Theta} |\bar{\ell}_n(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta})| \xrightarrow{p} 0 \quad \text{as } n \rightarrow \infty.$$

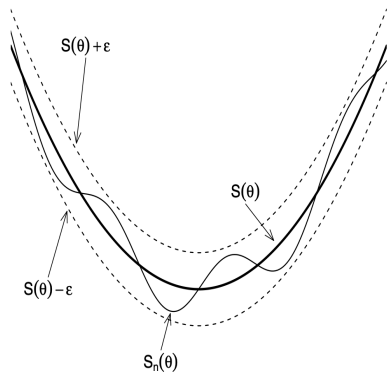
A sequence of vector random functions $\{\mathbf{h}_n(\cdot)\}$ converges uniformly in probability to a nonrandom function $\mathbf{h}(\cdot)$ if

$$\sup_{\boldsymbol{\theta} \in \Theta} \|\mathbf{h}_n(\boldsymbol{\theta}) - \mathbf{h}(\boldsymbol{\theta})\| \xrightarrow{p} 0 \quad \text{as } n \rightarrow \infty,$$

where $\|\cdot\|$ is the usual Euclidean norm.



(a) Non-Uniform Convergence



(b) Uniform Convergence

Existence and Consistency

定理 (Existence of MLE)

Suppose that (i) the parameter space Θ is a compact subset of \mathbb{R}^p , (ii) $\log f(x_i|\theta)$ is continuous in θ for each x_i , and (iii) $\log f(x_i|\theta)$ is a measurable function in x_i for all θ in Θ . Then there exists a measurable function $\hat{\theta}$ of the data that solves (1).

定理 (Consistency of MLE)

Let $\{X_i\}$ be i.i.d. Suppose that the conditions for Existence hold. Suppose, further, that

- (a) (*identification*) $\ell(\theta)$ is uniquely maximized on Θ at $\theta_0 \in \Theta$.
- (b) (*uniform convergence*) $\bar{\ell}_n(\theta)$ converges uniformly in probability to $\ell(\theta)$,

then $\hat{\theta} \xrightarrow{p} \theta_0$.

证明.

Let N be an open neighborhood in \mathbb{R}^p containing $\boldsymbol{\theta}_0$. Then $N^C \cap \Theta$ is compact. Therefore $\max_{\boldsymbol{\theta} \in N^C \cap \Theta} \ell(\boldsymbol{\theta})$ exists. Denote

$$\epsilon = \ell(\boldsymbol{\theta}_0) - \max_{\boldsymbol{\theta} \in N^C \cap \Theta} \ell(\boldsymbol{\theta}). \quad (2)$$

Let A_n be the event “ $|\bar{\ell}_n(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta})| < \epsilon/2$ for all $\boldsymbol{\theta}$.” Then

$$A_n \Rightarrow \ell(\hat{\boldsymbol{\theta}}) > \bar{\ell}_n(\hat{\boldsymbol{\theta}}) - \epsilon/2 \quad (3)$$

and

$$A_n \Rightarrow \bar{\ell}_n(\boldsymbol{\theta}_0) > \ell(\boldsymbol{\theta}_0) - \epsilon/2. \quad (4)$$

证明.

But, since $\bar{\ell}_n(\hat{\boldsymbol{\theta}}) \geq \bar{\ell}_n(\boldsymbol{\theta}_0)$, we have from (3)

$$A_n \Rightarrow \ell(\hat{\boldsymbol{\theta}}) > \bar{\ell}_n(\boldsymbol{\theta}_0) - \epsilon/2. \quad (5)$$

Therefore, adding both sides of (4) and (5), we obtain

$$A_n \Rightarrow \ell(\hat{\boldsymbol{\theta}}) > \ell(\boldsymbol{\theta}_0) - \epsilon. \quad (6)$$

Therefore, from (2) and (6) we can conclude

$$A_n \Rightarrow \hat{\boldsymbol{\theta}} \in N,$$

which implies

$$\text{Prob}(A_n) \leq \text{Prob}(\hat{\boldsymbol{\theta}} \in N).$$

But, since $\lim_{n \rightarrow \infty} \text{Prob}(A_n) = 1$ by (b), $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0$.



Asymptotic Normality

定理 (Asymptotic Normality of MLE)

Let $\{X_i\}$ be i.i.d. Suppose that the conditions for $\hat{\theta}$ to be consistent are satisfied. Suppose, further, that

- (1) θ_0 is in the interior of Θ ,
- (2) $f(x_i|\theta)$ is twice continuously differentiable in θ for any x_i ,
- (3) $\mathbb{E}[\frac{\partial}{\partial \theta} \log f(X_i|\theta_0)] = \mathbf{0}$ and $\mathcal{H}_{\theta} = \mathcal{I}_{\theta}$,
- (4) for some neighborhood \mathcal{N} of θ_0 , $\mathbb{E} \left[\sup_{\theta \in \mathcal{N}} \left\| \frac{\partial^2}{\partial \theta \partial \theta'} \log f(X_i|\theta) \right\| \right] < \infty$,
- (5) \mathcal{H}_{θ} is nonsingular.

Then $\hat{\theta}$ is asymptotically normal with

$$\text{Avar}(\hat{\theta}) = \mathcal{I}_{\theta}^{-1}.$$

证明.

By Mean Value Theorem, we have

$$\begin{aligned}
 \frac{\partial \bar{\ell}_n(\hat{\theta})}{\partial \theta} &= \frac{\partial \bar{\ell}_n(\theta_0)}{\partial \theta} + \frac{\partial^2 \bar{\ell}_n(\bar{\theta})}{\partial \theta \partial \theta'} (\hat{\theta} - \theta_0) \\
 &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i | \theta_0) + \left[\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta'} \log f(x_i | \bar{\theta}) \right] (\hat{\theta} - \theta_0) \\
 &= \mathbf{0},
 \end{aligned}$$

where $\bar{\theta}$ lies between θ_0 and $\hat{\theta}$. Therefore,

$$\sqrt{n}(\hat{\theta} - \theta_0) = - \left[\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta'} \log f(x_i | \bar{\theta}) \right]^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i | \theta_0).$$

□

证明.

By i.i.d. and assumption (4), we have

$$\begin{aligned}\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \log f(x_i | \boldsymbol{\theta}_0) &\xrightarrow{d} N(\mathbf{0}, \mathcal{I}_{\boldsymbol{\theta}}), \\ -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log f(x_i | \bar{\boldsymbol{\theta}}) &\xrightarrow{p} \mathcal{H}_{\boldsymbol{\theta}}.\end{aligned}$$

Therefore,

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} \mathcal{H}_{\boldsymbol{\theta}}^{-1} N(\mathbf{0}, \mathcal{I}_{\boldsymbol{\theta}}) = N(\mathbf{0}, \mathcal{I}_{\boldsymbol{\theta}}^{-1}).$$



$\text{Avar}(\hat{\boldsymbol{\theta}})$ can be consistently estimated by

- first estimator of $\text{Avar}(\hat{\boldsymbol{\theta}})$:

$$- \left\{ \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log f(x_i | \hat{\boldsymbol{\theta}}) \right\}^{-1}.$$

- second estimator of $\text{Avar}(\hat{\boldsymbol{\theta}})$:

$$\left\{ \frac{1}{n} \sum_{i=1}^n \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log f(x_i | \hat{\boldsymbol{\theta}}) \right] \left[\frac{\partial}{\partial \boldsymbol{\theta}'} \log f(x_i | \hat{\boldsymbol{\theta}}) \right] \right\}^{-1}.$$

Wald, LM, and LR Test

The Null Hypothesis

Let θ_0 be the p -dimensional model parameter. The null hypothesis can be expressed as

$$\mathbb{H}_0 : \underset{(r \times 1)}{\mathbf{c}(\theta_0)} = \mathbf{0}.$$

We assume that $\mathbf{c}(\cdot)$ is continuously differentiable. Also, let

$$\underset{(r \times p)}{\mathbf{C}(\theta)} \equiv \frac{\partial \mathbf{c}(\theta)}{\partial \theta'}.$$

We assume that

$$\underset{(r \times p)}{\mathbf{C}_0} \equiv \mathbf{C}(\theta_0) \text{ is of full row rank.}$$

Let $\hat{\theta}$ solve the unconstrained optimization problem. The constrained estimator, denoted $\tilde{\theta}$, solves

$$\max_{\theta \in \Theta} L_n(\theta) \quad \text{s.t.} \quad \mathbf{c}(\theta) = \mathbf{0}.$$

- Wald Statistic

$$W \equiv n\mathbf{c}(\hat{\boldsymbol{\theta}})' \left[\mathbf{C}(\hat{\boldsymbol{\theta}}) \widehat{\mathcal{I}}_{\boldsymbol{\theta}}^{-1} \mathbf{C}(\hat{\boldsymbol{\theta}})' \right]^{-1} \mathbf{c}(\hat{\boldsymbol{\theta}}) \xrightarrow{d} \chi^2(r)$$

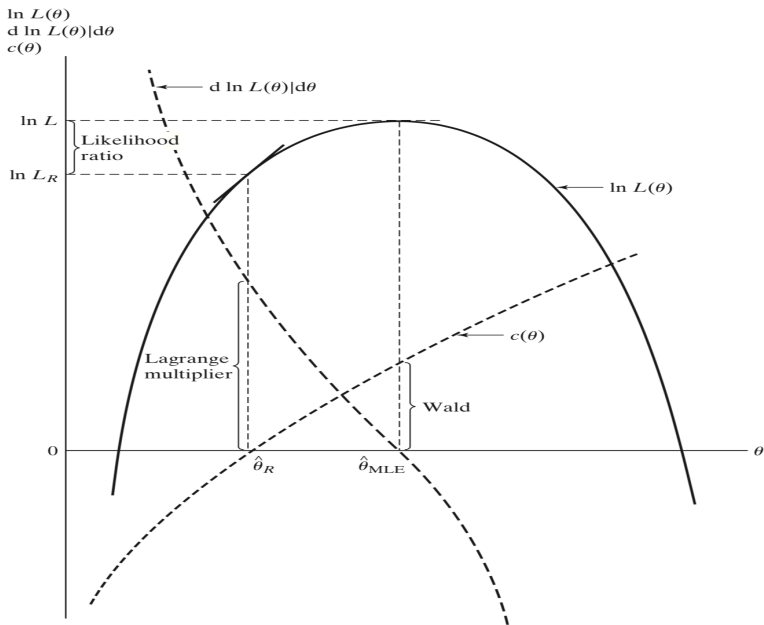
- Lagrange Multiplier (LM) Statistic (or, Score Test Statistic)

$$LM \equiv n \left(\frac{\partial \bar{\ell}_n(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} \right)' \widetilde{\mathcal{I}}_{\boldsymbol{\theta}}^{-1} \left(\frac{\partial \bar{\ell}_n(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} \right) \xrightarrow{d} \chi^2(r)$$

- Likelihood Ratio (LR) Statistic

$$LR \equiv 2n \cdot [\bar{\ell}_n(\hat{\boldsymbol{\theta}}) - \bar{\ell}_n(\tilde{\boldsymbol{\theta}})] \xrightarrow{d} \chi^2(r)$$

To summarize, the three statistics all converge in distribution to $\chi^2(r)$ under the null where r is the number of restrictions in the null.



References



Fumio Hayashi. (2000)

Econometrics, Chapter 1.

Princeton University Press, 2000.



Bruce E. Hansen. (2020)

Introduction to Econometrics, Chapter 10.

<https://www.ssc.wisc.edu/~bhansen/econometrics/>

▶ Link